# Exhibit M

### IN THE UNITED STATES DISTRICT COURT FOR THE EASTERN DISTRICT OF VIRGINIA ALEXANDRIA DIVISION

United States of America, et al.,

Plaintiffs,

v

Case No. 1:23-cv-00108 HON. LEONIE H. M. BRINKEMA

Google LLC,

Defendant.

EXPERT REPORT OF

GABRIEL WEINTRAUB, PH.D.

**DECEMBER 22, 2023** 

<sup>247</sup> Given the need for "sufficient" and the "right kind" of data for these algorithms, as outlined in this subsection, a company that falls out of the virtuous cycle will likely face difficulty in overcoming the hurdle when data scale, auction outcomes such as impressions won and revenue, and thickness are intertwined.

## III.C. Scale and Experimentation

100. Technology companies routinely use experiments to assess the impact of changing product (or service) features or introducing new features. For these experiments, they often use A/B tests, also known as online controlled experiments ("OCE").<sup>248</sup> An A/B test is a randomized experiment in which a fraction of subjects are randomly selected to receive a treatment while the remaining subjects form the control group that represents the baseline (which is often the status





(OpenX), October 26, 2023, 63:17–64:13 ("Q. Does OpenX vary its take rate by impression, meaning for a given publisher, OpenX might charge a different take rate for one impression than the other?...THE WITNESS: Yes. It varies – that varies by publisher in terms of how we apply our take rate, but a simple example is that we will often do something called fee squashing where we'll take a lower fee if we need to, to try and get an impression to run through our system. That would be an example....Q. How, if at all, does greater scale affect exchanges' ability to effectively set different take rates for different impressions? A....The greater the scale that you have, the more data you have. The more data you have, you can better decide what an appropriate fee would be on a given impression.").

Ron Kohavi, Diane Tang, and Ya Xu, *Trustworthy Online Controlled Experiments: A Practical Guide to A/B Testing*, (Cambridge University Press, 2020), 5 ("Controlled experiments have a long and fascinating history, which we share online (Kohavi, Tang and Xu 2019). They are sometimes called A/B tests, A/B/n tests (to emphasize multiple variants), field experiments, randomized controlled experiments, split tests, bucket tests, and flights. In this book, we use the terms *controlled experiments* and *A/B tests* interchangeably, regardless of the number of variants.") (emphasis in original).



101. A/B tests are considered to be a scientifically valid approach to determine the causal impact of a treatment.<sup>254</sup> These randomized OCE's are the "gold standard for establishing causality,"<sup>255</sup> which is why web-facing companies including eBay, Facebook, Google, and Microsoft (among



Ron Kohavi, Diane Tang, and Ya Xu, *Trustworthy Online Controlled Experiments: A Practical Guide to A/B Testing* (Cambridge: Cambridge University Press, 2020), 9 ("Randomized controlled experiments are the gold standard for establishing causality.").

Ron Kohavi, Diane Tang, and Ya Xu, Trustworthy Online Controlled Experiments: A Practical Guide to A/B Testing (Cambridge University Press, 2020), 9 ("Randomized controlled experiments are the gold standard for establishing causality.").

many others), use them to guide product development and accelerate innovation.<sup>256</sup> When designing experiments, there are best practices to ensure that the results are statistically valid;



Ron Kohavi, Alex Deng, Brian Frasca, Toby Walker, Ya Xu, and Nils Pohlmann, "Online Controlled Experiments at Large Scale," *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (2013): 1–9, at 1 ("Many web-facing companies use online controlled experiments to guide product development and prioritize ideas, including Amazon [1], eBay, Etsy [2], Facebook, Google [3], Groupon, Intuit [4], LinkedIn, Microsoft [5], Netflix [6], Shop Direct [7], StumbleUpon [8], Yahoo, and Zynga [9].").





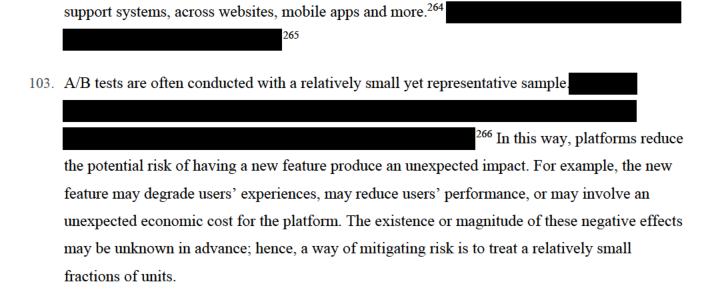
102. In this way, A/B tests have become an integral part of technological companies' innovation engines. These companies rely on experiments to gauge the effects of potential changes to their platforms before implementation. Organizations use these experiments as a key tool for developing frameworks and testing ideas to understand the impact of the services they provide; they have become a standard component of managerial decision making. Companies like eBay, Facebook, Microsoft and Google run thousands to tens of thousands of experiments every year, testing everything from changes to the user interface, to customer



Michael Luca, and Max H. Bazerman, *The Power of Experiments: Decision Making in a Data-Driven World* (Cambridge: MIT Press, 2021): vii-viii ("These days, companies like Google wouldn't dare make a major change in their platforms without first looking at experiments to understand how it would influence user behavior.").

Michael Luca and Max H. Bazerman, The Power of Experiments: Decision Making in a Data-Driven World (Cambridge: MIT Press, 2021), vii–viii ("From startups to international conglomerates to government agencies, organizations have a new tool to develop frameworks and test ideas, and to understand the impact of the products and services they are providing.").

Michael Luca and Max H. Bazerman, *The Power of Experiments: Decision Making in a Data-Driven World* (Cambridge: MIT Press, 2021), 62 ("But perhaps no sector has embraced the experimental method more than the tech sector, where it is now a standard component of managerial decision making.").



104. Statistical power is an important criterion to consider when conducting statistically sound experiments. <sup>267</sup> Suppose an experiment is conducted to detect an X percent (or larger) improvement in an outcome of interest caused by a new feature. This "X percent" could be the smallest detectable effect that the company would use as a threshold to decide whether to launch the feature system-wide. Statistical power refers to the probability of detecting such an effect

265

Ron Kohavi, Diane Tang, and Ya Xu, Trustworthy Online Controlled Experiments: A Practical Guide to A/B Testing (Cambridge University Press, 2020), 5 ("Online controlled experiments are used heavily at companies like Airbnb, Amazon, Booking.com, eBay, Facebook, Google, LinkedIn, Lyft, Microsoft, Netflix, Twitter, Uber, Yahoo!/Oath, and Yandex (Gupta et al. 2019). These companies run thousands to tens of thousands of experiments every year, sometimes involving millions of users and testing everything, including changes to the user interface (UI), relevance algorithms (search, ads, personalization, recommendations, and so on), latency/performance, content management systems, customer support systems, and more. Experiments are run on multiple channels: websites, desktop applications, mobile applications, and e-mail.").

Ron Kohavi, Diane Tang, and Ya Xu, Trustworthy Online Controlled Experiments: A Practical Guide to A/B Testing, (Cambridge University Press, 2020): 171 ("In practice, it is common that an experiment goes through a ramping process to control unknown risks associated with new feature launches (aka. controlled exposure). For example, a new feature may start by exposing the Treatment to only a small percentage of users. If the metrics look reasonable and the system scales well, then we can expose more and more users to the Treatment. We ramp the traffic until the Treatment reaches desired exposure level.")

Ron Kohavi, Diane Tang, and Ya Xu, Trustworthy Online Controlled Experiments: A Practical Guide to A/B Testing (Cambridge University Press, 2020), 30 ("Practically speaking, you want enough power in your experiment to be able to conclude with high probability whether your experiment has resulted in a change bigger than what you care about.").

when in fact one exists.<sup>268</sup> There are multiple factors that affect the statistical power of an experiment, such as the number of observations, where fewer observations lead to lower statistical power, and variance in the outcome of interest, where higher variance leads to lower statistical power.<sup>269</sup>

- 105. A significant challenge in running A/B tests for technology platforms that run many experiments simultaneously is to secure enough users to achieve a desired level of statistical power.

  Technology companies are constantly launching, changing, and iterating on their products.<sup>270</sup>

  Each individual change may have a small effect, but the compounding effect of many small changes can be significant.<sup>271</sup> Therefore, detecting small effects reliably (or in a statistically significant way) is critical to be able to continuously experiment and improve.
- 106. Statistical theory shows the size of the experimental sample should become larger as the level of the effect that an experimenter wants to detect becomes smaller.<sup>272</sup> For example, if a company

- Ron Kohavi, Diane Tang, and Ya Xu, *Trustworthy Online Controlled Experiments: A Practical Guide to A/B Testing* (Cambridge University Press, 2020), 30 ("Usually, we get more power when the sample size is larger."), and 189–190 ("Assuming Treatment and Control are of equal size, the total number of samples you need to achieve 80% power can be derived from the power formula above, and is approximately as showing in Equation 17.8 (van Belle 2008):  $n \approx \frac{16\sigma^2}{\delta^2}$  where,  $\sigma^2$  is the sample variance, and  $\delta$  is the difference between Treatment and Control.").
- Ron Kohavi, Diane Tang, and Ya Xu, Trustworthy Online Controlled Experiments: A Practical Guide to A/B Testing (Cambridge University Press, 2020), 13 ("Features are built because teams believe they are useful, yet in many domains most ideas fail to improve key metrics. Only one third of the ideas tested at Microsoft improved the metric(s) they were designed to improve (Kohavi, Crook and Longbotham 2009). Success is even harder to find in well-optimized domains like Bing and Google, whereby some measures' success rate is about 10-20% (Manzi 2012).").
- Ron Kohavi, Diane Tang, and Ya Xu, Trustworthy Online Controlled Experiments: A Practical Guide to A/B Testing (Cambridge University Press, 2020), 14 ("[Google] ran hundreds of controlled experiments and multiple iterations; some across all markets, and some long term in specific markets to understand the impact on advertisers in more depth. This large backend change and running controlled experiments ultimately validated how planning multiple changes and layering them together improved the user's experience by providing higher quality ads, and improved their advertiser's experience moving towards lower average prices for the higher quality ads.").
- Ron Kohavi, Alex Deng, Brian Frasca, Toby Walker, Ya Xu, and Nils Pohlmann, "Online Controlled Experiments at Large Scale," *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (2013): 1–9, at 1 ("While running online controlled experiments require a sufficient number of users, teams working on products with thousands to tens of thousands of users (our general guidance is at least thousands of active users) are typically looking for larger effects, which are easier to detect than the

<sup>&</sup>lt;sup>268</sup> Ron Kohavi, Diane Tang, and Ya Xu, *Trustworthy Online Controlled Experiments: A Practical Guide to A/B Testing* (Cambridge University Press, 2020), 30 ("*Statistical power* is the probability of detecting a meaningful difference between the variants when there really is one (statistically, reject the null when there is a difference).") (emphasis in original).

does not want to introduce a new feature unless an expected effect has a 1 percent lift on traffic, it should run a larger sample size experiment than if their target effect is 2 percent. As a result, large sample sizes are required for conducting reliable A/B tests for individual changes with relatively small effects.<sup>273</sup>

107. For this reason, it can be challenging for firms with limited scale to run informative experiments. 274 For technology platforms, the unit of an observation in an experiment typically corresponds to a user using the platform (e.g., a bidder), or to an event (e.g., a query, an auction or an impression). Hence, typically platforms with smaller scale need to run experiments over a longer time period before the number of users or events included in the experiment reaches the sample size required to detect the desired effect. 275 Running such long experiments can be impractical when a firm's objective is to iterate many small changes quickly. For example, most large tech companies run a single experiment for two weeks or less. 276 Consider another smaller company that has 20 percent of the traffic of the large tech company and uses the same treatment allocation fraction as the large company. Then, everything else equal, this smaller company

small effects that large sites worry about. For example, to increase the experiment sensitivity (detectable effect size) by a factor of 10, say from 5% to 0.5%, you need  $10^2 = 100$  times more users.").

see also, "Find top products in Learning Management Systems (LMS) category," LinkedIn, accessed December 20, 2023

Learning Management Systems (LMS) category," LinkedIn, accessed December 20, 2023, https://www.linkedin.com/products/categories/learning-management-systems.

Ron Kohavi, Diane Tang, and Ya Xu, *Trustworthy Online Controlled Experiments: A Practical Guide to A/B Testing* (Cambridge University Press, 2020), 33 ("**More users:** In the online experiments, because users trickle into experiments over time, the longer the experiment runs, the more users the experiment gets. This usually results in increased statistical power") (emphasis in original);

Somit Gupta, Ronny Kohavi, Diane Tang, Ya Xu, et al., "Top Challenges from the First Practical Online Controlled Experiments Summit," ACM SIGKDD Explorations Newsletter 21, no. 1 (2019): 20–35, at 21–22 ("There are many interesting and challenging open questions for OCE results analysis. While most experiments in the industry run for 2 weeks or less, we are really interested in detecting the long-term effect of a change... At Microsoft, while most experiments do not run for more than two weeks, it is recommended to run an experiment longer if novelty effects are suspected and use data from the last week to estimate the long-term treatment effect...At Twitter, a similar practice is followed. An experiment at Twitter may run for 4 weeks and data fromlast two weeks is analyzed.").

<sup>273</sup> Ron Kohavi, Diane Tang, and Ya Xu, Trustworthy Online Controlled Experiments: A Practical Guide to A/B Testing (Cambridge University Press, 2020), 30 ("Usually, we get more power when the sample size is larger.").

Wigh

Gabriel Weintraub, Ph.D.

December 22, 2023

#### Expert Report of Gabriel Weintraub (December 22, 2023)--Errata

Page	Paragraph	Footnote Original	Corrected	Reason
<u>7</u>	16	6 "ad.spend"	"ad spend"	Туро
11	22	10 "See, e.g., "Internet Advertising Revenue Report," PwC, iab, April 2023,	"See, e.g., "Internet Advertising Revenue Report," PwC, iab, April 2023,	Clarification
		https://www.iab.com/wpcontent/	https://www.iab.com/wpcontent/	
			p uploads/2023/04/IAB_PwC_Internet_Advertising_Revenue_Report_2022. ا	0
		df"	df, at 15"	
12	25	. "The tools to run RTB auctions for the purchase and sale of online display	"The tools to run RTB auctions for the purchase and sale of online display	Туро
		ads <b>is</b> called"	ads <b>are</b> called"	
17	33	41 "(e.g., the average spend for each conversion should not exceed a pre-	"(e.g., the average spend for each conversion should not exceed a pre-	Clarification
		specified target).")."	specified target).") (emphasis in original)."	
19	34	47 "at 39 <b>68</b> "	"at 39 <b>52</b> "	Correction
23	43	71 "using machine algorithms"	"using machine learning algorithms"	Clarification
24	44	77 "no. 1 (2014)"	"no. 1 (2015)"	Туро
25	46	88 "EC '22: Proceedings of the <b>2023</b> ACM Conference on Economics and	"EC '22: Proceedings of the <b>23rd</b> ACM Conference on Economics and	Туро
25	4.0	Computation <b>2023</b> "	Computation"	Tuno
25 25	46 46	88 "no. 1 (2023)" 90 "EC '22: Proceedings of the 2023 ACM Conference on Economics and	"no. 1 (2022)"  "EC '22: Proceedings of the 23rd ACM Conference on Economics and	Туро
25	46	Computation <b>2023</b> "	Computation"	туро
25	46	90 "no. 1 ( <b>2023</b> )"	"no. 1 (2022)"	Туро
26	46	90 "(2023): 1–35, at 6"	"(2023): 1–35, at 6– <b>7</b> "	Clarification
26	46	91 "EC '22: Proceedings of the <b>2023</b> ACM Conference on Economics and	"EC '22: Proceedings of the <b>23rd</b> ACM Conference on Economics and	Туро
20	40	Computation 2023"	Computation"	Туро
26	46	91 "no. 1 ( <b>2023</b> )"	"no. 1 (2022)"	Туро
		32 No. 2 (2020)		1,700
28	51	. "Google Ad Manager ("GAM"), which includes a publisher ad server, and	"Google Ad Manager ("GAM"), which includes a publisher ad server and	Туро
		an ad exchange/SSP."	an ad exchange/SSP."	71-
32	59	132 " <i>6</i> <sup>1</sup> (x) "	" <i>β</i> '( <i>x</i> ) "	Туро
36	68	143 "("we use the market thickness to represent the average number of ads	"("we use the market thickness to represent the average number of ads	Correction
		competing for user impressions on an online advertising platform.")	competing for user impressions on an online advertising platform.");"	
_		(emphasis in original);"		
38	70	154 " <b>955</b> –1025"	" <b>965</b> –1025"	Туро
38	70	157 " <b>955</b> –1025"	" <b>965</b> –1025"	Туро
56	93	225 "1849–1864, at <b>1850</b> "	"1849–1864, at <b>1851</b> "	Туро
60	98	243 "209–218, at <b>217</b> "	"209–218, at <b>209</b> "	Correction
60	98	244 "209–218, at <b>217</b> "	"209–218, at <b>216</b> "	Correction
	460	264 II/C I - I AUT D 2024)II	III/Co. I i I o MIT Dono 2000\III	
64	102	261 "(Cambridge: MIT Press, <b>2021</b> )"	"(Cambridge: MIT Press, 2020)"	Туро
64	102	262 "(Cambridge: MIT Press, <b>2021</b> )"	"(Cambridge: MIT Press, 2020)"	Туро
64	102	263 "(Cambridge: MIT Press, <b>2021</b> )"	"(Cambridge: MIT Press, <b>2020</b> )"	Typo
66	104	269 "and 189 <b>–190</b> "	"and 189"	Clarification

#### Expert Report of Gabriel Weintraub (December 22, 2023)--Errata

Page	Paragraph	Footnote Original	Corrected	Reason
68	109	279 "at 3933"	"at 3933 <b>–3934</b> "	Clarification
68	109	279 "(2017): 500–522."	"(2017): 500–522 <b>, at 500, 510</b> ."	Clarification
70	112	284 " <b>23:6</b> –24:6"	" <b>24:1</b> –24:6"	Clarification
70	112	. "John Gentry of OpenX named personnel costs alongside infrastructure costs, and real estate costs as costs that need <b>to covered</b> in order to have the "ability to invest and grow the business.""	"John Gentry of OpenX named personnel costs alongside infrastructure costs and real estate costs as costs that need <b>to be covered</b> in order to have the "ability to invest and grow the business.""	Туро
				al if: ii
73	115	296 "247:17–248: <b>19</b> "	"247:17–248: <b>7</b> "	Clarification
- 11				
	_			
115	169	39 <u>6</u> "189:5–190: <b>18</b> "	"189:5–190: <b>3</b> "	Clarification
_				
	_			
146	220	. "Recall that in first price auctions, it is optimal for bidders to reduce their bids (Section <b>II.A.1</b> )."	"Recall that in first price auctions, it is optimal for bidders to reduce their bids (Section II.C.1)."	Туро
153	220	FF2 "Denosition of John Control On-1911 Oct-1-1-26 2022 24:45 24 22 1//2	Switch with tout in factors FF2	Connoction
152	229	552 "Deposition of John Gentry (OpenX), October 26, 2023, 21:15–21:22 ("Q. How, if at all, did the decrease in spending by DV360 on OpenX affect OpenX?THE WITNESS: It was a devastating impact to the company, resulting in severe financial consequences. We had to execute a large layoff in December of 2018 and had a lot of negative effects as a result of	Switch with text in footnote 553.	Correction

#### Expert Report of Gabriel Weintraub (December 22, 2023)--Errata

Page	Paragraph	Footnote Original	Corrected	Reason
152	229	"Deposition of John Gentry (OpenX), October 26, 2023, 22:17–23 ("Our revenues went from - we were about in 2017; 2018 we were about because the first half of the year was strong; in 2019, we were down to about; and in 2020, the combination of difficulties we had had plus COVID took us down to about in net revenue in 2020.")."	Switch with text in footnote 552.	Correction
152	229	. "The CEO of OpenX, John Gentry, testified that OpenX's revenues fell from in 2017 to 2019."	"The CEO of OpenX, John Gentry, testified that OpenX's revenues fell from in 2017 to 2019."	Туро
C-1	3.a.ii.	3 "Letter from David R. Pearl to <b>Kelly Garcia</b> , September 8, 2023, 2."	"Letter from David Pearl to <b>Michael Freeman</b> , September 8, 2023, 2."	Correction
D-10	22	17 "Letter from Julie Elmer to John Hogan, August 19, 2022, 5–6   8 ("	"Letter from Julie Elmer to John Hogan, August 19, 2022, 5–6 ("	Туро
E-2	2	. "exchanges' bids have the same variance that is characterized by the spread parameter"	"exchanges' bids have variance that is characterized by the same spread parameter"	Clarification
E-21	42	21 "Letter from David Pearl to Michael J. Freeman, July 28, 2023 ("	"Letter from David Pearl to Michael J. Freeman, July 28, 2023, 2 ("	Clarification
F-2	3	5 "John Rice, "11.2 Comparing Two Independent Samples", in Mathematical Statistics and Data Analysis, 3rd ed. (Duxbury: Thomson Brooks/Cole, 2007)."	"John Rice, "11.2 Comparing Two Independent Samples", in Mathematical Statistics and Data Analysis, 3rd ed. (Duxbury: Thomson Brooks/Cole, 2007), 421–444."	Clarification
F-5	11	12 ("Rule of Thumb[:] The basic formula is n=16/ $\Delta$ ^2 (2.3).").	("Rule of Thumb[:] The basic formula is n=16/ $\Delta$ ^2 (2.3) where $\Delta$ =( $\mu$ _0- $\mu$ _1)/ $\sigma$ = $\delta$ / $\sigma$ (2.4).").	Clarification
F-5	11	12 N/2=16 ( p_p (1-p_p))/(p_1-p_0 )	N/2=16 ( p_p (1-p_p))/(p_1-p_0 ) <b>^2</b>	Correction

January 13, 2024

# Case 1:23-cv-00108-LMB-JFA Document 640-13 Filed 05/17/24 Page 14 of 14 PageID# HIGHLY CONFIDENTIAL - 1/20/01/25 TO PROTECTIVE ORDER

#### **Expert Report of Gabriel Weintraub (December 22, 2023)--Supplemental Errata**

Page Paragraph	Footnote Original	Corrected	Reason

February 23, 2024